

OUTLIER DETECTION IN TOTAL PHOSPHORUS CONCENTRATION  
DATA FROM SOUTH FLORIDA RAINFALL

HOSUNG AHN

WRF-359

*Made in United States of America*  
Reprinted from JOURNAL OF THE AMERICAN WATER RESOURCES ASSOCIATION  
Vol. 35, No. 2, April 1999  
Copyright © 1999 by the American Water Resources Association



OUTLIER DETECTION IN TOTAL PHOSPHORUS CONCENTRATION  
DATA FROM SOUTH FLORIDA RAINFALL<sup>1</sup>Hosung Ahn<sup>2</sup>

**ABSTRACT:** Atmospheric deposition can be a significant source of phosphorus to South Florida's aquatic system. Deposition samples are often contaminated to varying degrees by bird droppings or other foreign materials. This study attempted to use statistical and other methods to detect and remove the outliers in the rain-borne total phosphorus concentration data. Some outliers in the data were identified using field notes derived from visual inspection of the samples. Outlier detection statistics based on a simple linear regression were then used for additional data screening. As a result of these analyses, about 35 percent of the observed values were identified as outlying data which needed to be removed prior to further data analyses. Based on detected outliers in the data from 15 monitoring sites, a lumped cutoff value of 130  $\mu\text{g/L}$  was determined. This lumped cutoff value may be useful for further quality control and analyses of the data from the region.

**(KEY TERMS:** total phosphorus concentration; atmospheric deposition; sample contamination; outlier detection; linear regression.)

## INTRODUCTION

The management of phosphorus inputs to the South Florida ecosystems has become an increasing concern resulting in the need for accurate monitoring and analysis of phosphorus distributions. Atmospheric deposition can be a significant source of phosphorus to ecosystems in South Florida, where most water bodies are large and shallow. Atmospheric deposition is commonly sampled as wet and dry forms separately. Wet deposition is from rain, while dry deposition occurs as dustfall under dry conditions.

The South Florida Water Management District (District) has been collecting atmospheric deposition data in the region since the early 1970s. The monitoring program was significantly improved in 1992 by deploying wet/dry collectors (Aerochem Metrics Model

301) and adopting a standard operating procedure for data collection and processing in accordance with recommendations of the National Atmospheric Deposition Program (NADP) (Bigelow, 1984; Bigelow and Dossett, 1988). Currently, there are 19 atmospheric deposition monitoring sites operated by the District. Both wet and dry deposition samples have been collected at weekly intervals and analyzed at the District's laboratory in order to determine the level of nutrients and major ions.

Because most monitoring sites are located at or near marshes, contamination of the samples by bird droppings, insects, and debris is very common and problematic. This type of contamination results in high total phosphorus (TP) concentrations and adversely affects computation of the summary statistics of the data. Improvements in sample processing and installation of bird deterrents (Asman *et al.*, 1982; van Wyk and Stock, 1991) have lowered the frequency of contamination, but have not eliminated the problem completely.

The purpose of this paper is to present a two-step approach used to identify outliers in wet TP concentration data in rainfall samples. The first step employed to detect outliers was an examination of field notes, especially the visual descriptions of the samples during collection and analysis. The field note information provided a binary decision of whether each sample was contaminated or not based on the type of data flags. The contaminated data identified through this step were removed prior to further analysis. Because abnormally high wet TP concentrations were found among the remaining data, the second step involved an attempt to detect outliers using outlier detection statistics. A statistical outlier detection

<sup>1</sup>Paper No. 98032 of the *Journal of the American Water Resources Association*. Discussions are open until December 1, 1999.

<sup>2</sup>Lead Hydrologist, Resources Assessment Division, WRE, South Florida Water Management District, 3301 Gun Club Road, MS 7120, West Palm Beach, Florida 33406 (E-Mail: hosung.ahn@sfwmd.gov).

method is described and then applied to the wet TP data collected from the District's atmospheric deposition monitoring sites.

## DEFINITION AND METHOD

Outliers are data points that appear to deviate markedly from other members of the sample group in which they occur (Grubbs, 1969; Beckman and Cook, 1983; Barnett and Lewis, 1984). In relation to statistical analyses, Rousseeuw and van Zomeren (1990) defined outliers as observations that deviate from the estimates by a statistical model suggested by most of a data set, which is a mixture of clean and contaminated data. The latter definition implies that, to detect outliers, a statistical model can be used to define the differences (residuals) between observations and estimates, and the residuals can then be used as indicators of aberrant data.

There are a variety of statistical methods for detecting outliers (Barnett and Lewis, 1984; Beckman and Cook, 1983). One way of detecting outliers is to set an outlier bound at either two or three standard deviations from the mean. However, this simple method cannot be used here because the prior population statistics of uncontaminated data are unknown. Statistical modeling methods for detecting outliers rely on sample statistics. These methods include linear regression, multivariate (Rousseeuw and Van Zomeren, 1990; Atkinson and Mulira, 1993; Hadi, 1994; others), and time series analysis (Beckman and Cook, 1983; Chib and Tiwari, 1994; Tiwari and Dienes, 1994). However, the multivariate method is not applicable here because many wet TP concentration data are randomly missing so that the number of complete data (data having no missing at each time step) becomes quite small. Moreover, preliminary analysis revealed a weak serial correlation in these data sets, making the time series analysis also inappropriate for the data. Thus, linear regression methods were considered here.

### *Detecting Multiple Outliers with Linear Regression*

The linear regression method detects outliers by forming a clean subset (which is a presumed subset having no outlier in it), fitting the regression for the clean subset, and testing for outliers relative to the clean subset based on a test statistic. The clean subset should produce, among all possible subsets, the smallest residual sum of squares. Finding a clean subset from a given data set is not easy. That is, to find the clean subset having a size of  $i$  from a sample

data set having a size of  $n$ , it is necessary to fit a regression model to each of the  $\binom{n}{i}$  possible subsets where  $\binom{n}{i}$  is the number of combinations. Finding the minimum residual sum of squares requires extensive computations that may not even be feasible, especially for a large  $n$ .

Three approaches are available for finding a clean subset more effectively: a random search algorithm (Rousseeuw and van Zomeren, 1990), a forward search algorithm (Hadi and Simonoff, 1993; Atkinson, 1994) and an elemental set algorithm (Hawkins and Simonoff, 1993). The forward search algorithm suggested by Hadi and Simonoff (HS) (1993) was used here, because it is relatively simple and efficient computationally (Woodruff and Roche, 1993; Atkinson, 1994). The HS method based on a forward search algorithm starts by finding an initial clean subset  $M$  of size  $h = (n + k - 1)/2$ , then it searches a clean subset iteratively with increasing the size of  $M$  with checking outliers by  $t$ -statistics (refer to Figure 1 for details).

In the linear regression method, the Studentized residuals are often used to test multiple outliers (Beckman and Cook, 1983). To introduce the Studentized residual in the HS method, let us assume that a data set having a size of  $n$  is fitted by a simple regression model:

$$Y = X\beta + e \quad (1)$$

where  $Y$  is an  $n$ -vector of responses,  $X$  is an  $(n \times k)$  matrix representing  $k$  explanatory (independent) variables with a rank of  $k < n$ ,  $\beta$  is a  $k$ -vector of regression parameters, and residuals  $e$  is an  $n$ -vector of Gaussian random noises with  $N(0, \sigma_e^2)$ . In an atmospheric deposition context,  $Y$  could represent a set of wet TP measurements at a given site, and  $X$  may be a set of concurrent measurements from nearby sites. The least square estimates of  $\beta$  and  $\sigma_e^2$  are given by  $\hat{\beta} = (X^T X)^{-1} X^T Y$  and  $\hat{\sigma}_e^2 = e^T e / (n - k)$ , respectively, where  $( )^T$  denotes a matrix transpose.

A clean subset  $M$  is assumed where  $X_M$  and  $Y_M$  are the components in  $M$ , and  $\beta_M$  and  $\sigma_M^2$  are the corresponding regression parameters and residual variance, respectively. A Studentized residual  $d_i (i = k, \dots, n)$  for  $M$  is then defined (Hadi and Son, 1990; Hadi, 1992) as

$$d_i = D_i / \sigma_M = |y_i - x_i^T \hat{\beta}_M| / \left[ \sigma_M \sqrt{1 - x_i^T (X_M^T X_M)^{-1} x_i} \right], \quad \text{if } i \in M,$$

$$- |y_i - x_i^T \hat{\beta}_M| / \left[ \sigma_M \sqrt{1 + x_i^T (X_M^T X_M)^{-1} x_i} \right], \quad \text{if } i \notin M, \quad (2)$$

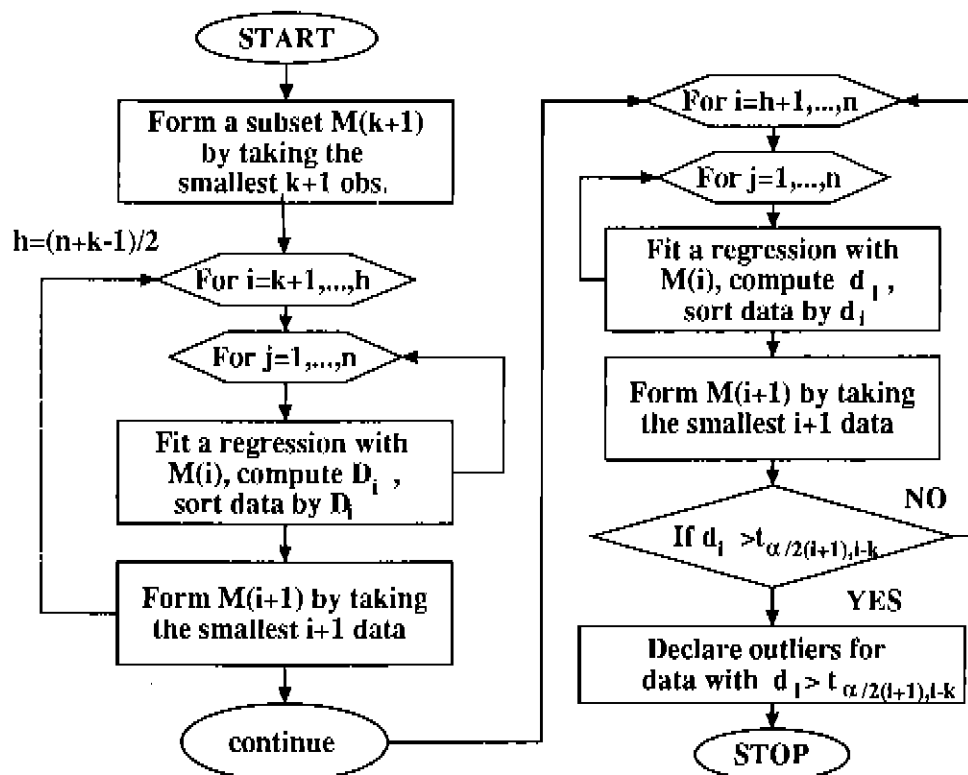


Figure 1. Forward Search Algorithm Proposed by Hadi and Simonoff (1993).

In particular, residuals ( $d_i$ ) for  $i \in M$  cases are the scaled prediction error relative to the subset  $M$ . Because the residuals follow a Student  $t$ -distribution, outliers in  $Y$  are tested with the statistics  $t_{\alpha/2(i+1),i-k}$  where  $\alpha/2(i+1)$  is the probability level and  $(i-k)$  is the degree of freedom of Student's  $t$ -distribution. All observations where  $d_i \geq t_{\alpha/2(i+1),i-k}$  are considered outliers.

The significance level,  $\alpha$ , is the only constant that needs to be determined before analysis. For outlier detection, commonly used significance levels are 10 percent (Cook, 1977; Jain, 1981), 5 percent (Hadi and Simonoff, 1993; Atkinson, 1994), 2.5 percent (Rousseeuw and van Zomeren, 1990), and 1 percent (Jain, 1981). However, the result of outlier detection is not sensitive to  $\alpha$  as will be shown later.

## WET TP CONCENTRATION DATA

Among 19 atmospheric deposition monitoring sites operated by the District (Figure 2), only 15 sites were analyzed in this investigation because the data from the remaining four sites (ENR101, ENR203, ENR301, ENR401) have relatively high rates of contamination: about 70 per cents of wet TP data from these four

sites are greater than 130  $\mu\text{g/L}$ . The maximum period of record for the 15 sites ranges from April 7, 1992, to October 22, 1996, but the actual record lengths vary from site to site owing to periodic expansion of the monitoring program: five sites (ENR, OKEEFS, S-140, S-65A, S-7) started sampling in April 1992, six sites (BG1, BG2, ENPRC, S-127, S-131, S-310) in September 1993, one site (S-308) in August 1994, and the remaining three sites (G-36, L-6, L-67A) in August 1995.

Figure 3 summarizes a schematic of the data classification and analyzing processes with the corresponding number of data in each class, where the numbers of data after Step I and II are from the result of the following two sections. The wet TP concentration data have been collected only for rainy weeks, but some of these data are missing (no observation) due to instrumental failures or other reasons.

## STEP I: ANALYZING THE DATA FLAGS

According to the District's standard operating procedure for handling wet deposition samples, all insects and animals are eliminated from the samples in the field while all other non-representative matters

## Area Map

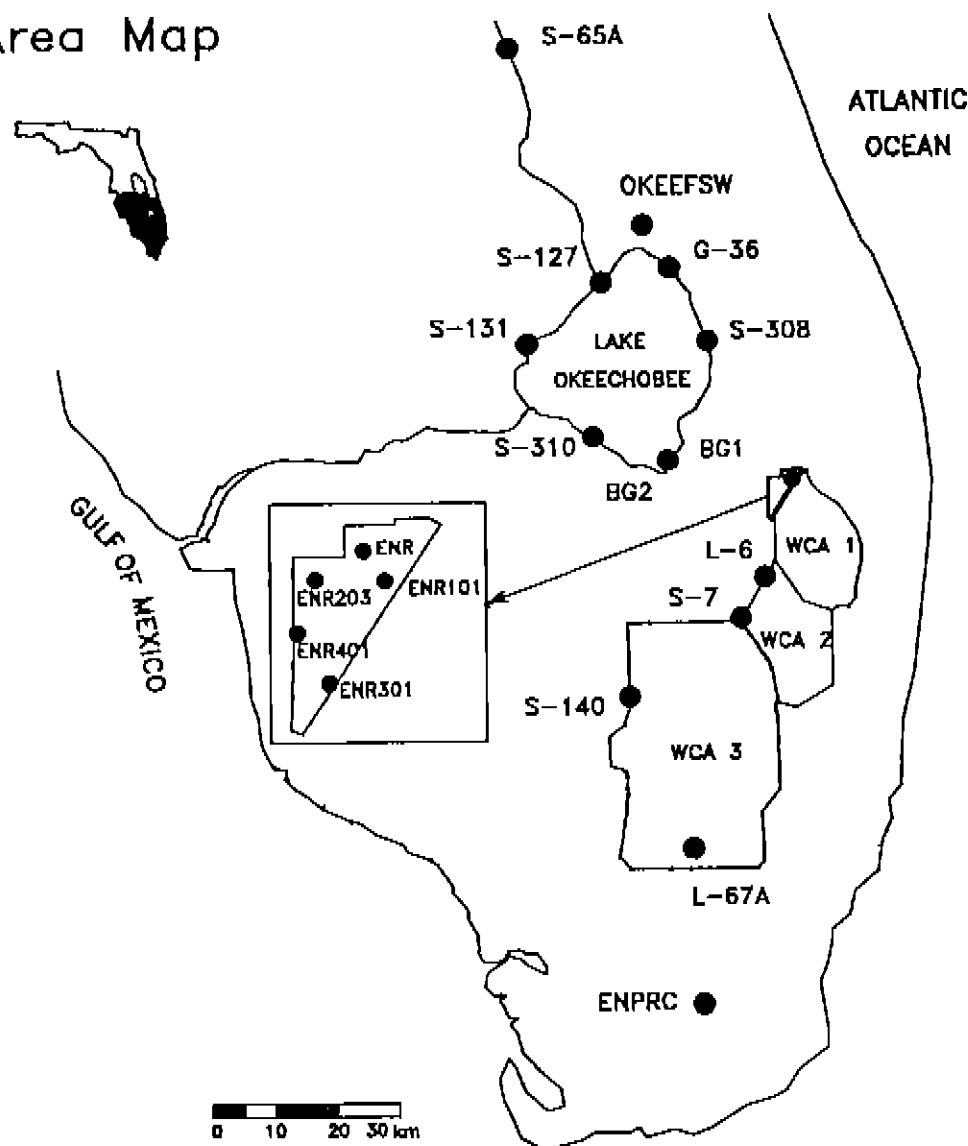


Figure 2. Location Map Showing the Atmospheric Deposition Monitoring Sites Operated by the South Florida Water Management District, Where WCA Stands for Water Conservation Area.

This monitoring network has been operated independently from the NADP network.

are removed at the laboratory with teflon tweezers (SFWMD, 1996). Visual contaminants are recorded in field notes and permanently logged into a computer data base to put flags on the measured wet TP concentration values.

To identify outlying data, an attempt was made to use the data flag information to salvage the contaminated data by separating the true atmospheric deposition component and contamination components. That is, the wet TP data were sorted according to the type of flags (possible contamination sources) and the mean ( $\bar{x}$ ) and variance ( $s^2$ ) of the TP values for each flag type were computed. The resulting statistics are

not presented here, but the result of this analysis revealed that a firm relationship between TP value and the type of flags was impossible to establish because of high variability (several orders of magnitude) of contaminated TP values and multiple flags in a TP value. Regarding the multiple flags, there is a total of 1585 cases of noted flags in 933 data among 2460 data points from 15 sites.

However, it was observed from the above analysis that some contamination sources were associated with high TP values more consistently than others. Thus, the contamination sources were grouped into six distinct categories which in turn were classified as

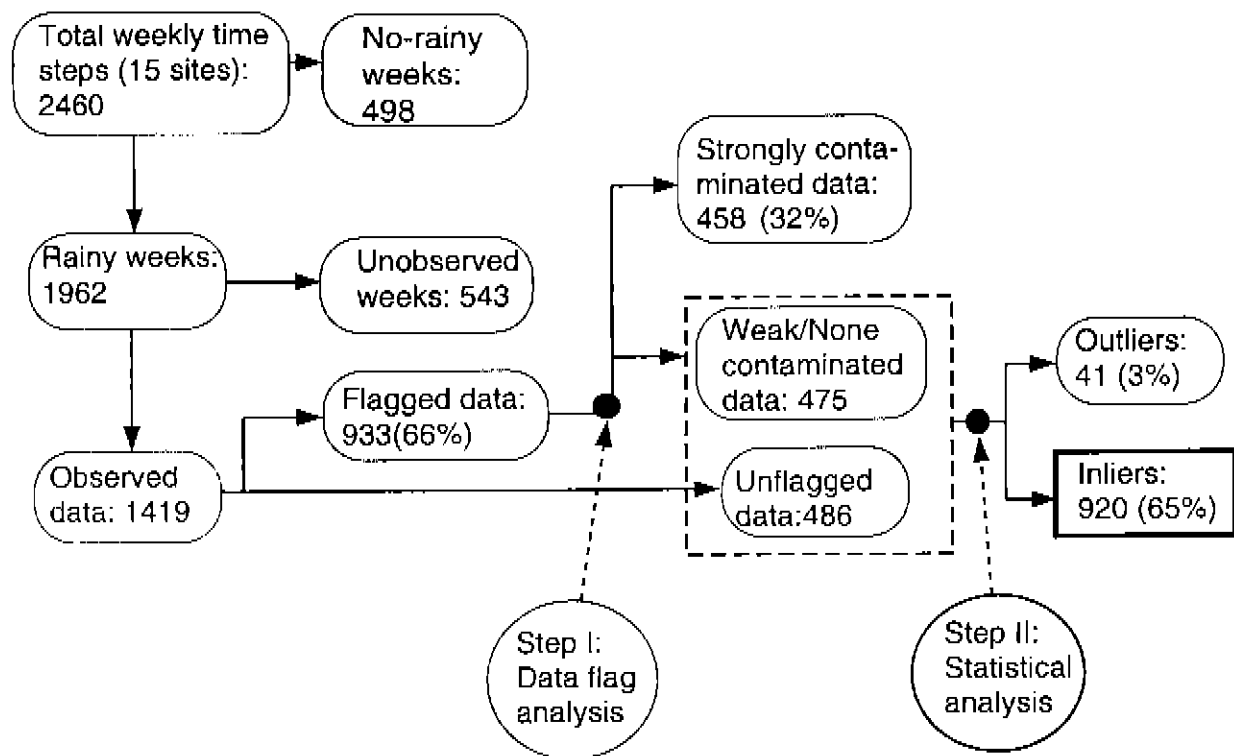


Figure 3. Classification of Wet TP Concentration Data, Along With the Corresponding Number of Samples and the Ratio (percent) of Classified Samples to Total Observed Samples in Parenthesis.

TABLE 1. Classification of Contamination Sources in Wet TP Concentration Data.  
+S indicates a strong contaminant and WN is a weak or no contaminant.

Class	Contamination sources	Type <sup>+</sup>	Percent
1. Bird-Dropping	Bird droppings, decomposed feces, urine, feces, organic materials, stains, suspended solution of feces	S	12
2. Dirt	Dirts, dust and dirt	S	27
3. Insect and Animal(I)	Insect body parts, frogs	S	4
4. Insect and Animal(II)	Ants, bees, beetles, cobwebs, feathers, flies, gnats, live insects, lizard, mosquitos, spiders, wasps, wings, other insects (heavy, moderate, light)	WN	28
5. Vegetation	Algae, berries, cut grass, vegetation	WN	5
6. Miscellaneous	Ash, dust, pollen, condensation, dew, unidentified	WN	24

either a strong contaminant (S) or weak/none contaminant (W) as given in Table 1. For instance, bird droppings were the most strong contaminant ( $\bar{x} \pm s = 1940 \pm 4700 \mu\text{g/L}$ ). On the other hand, the contribution from live insects or vegetation to TP values was much smaller. The presence of cobwebs associated with stains and insect body parts often resulted in high TP values ( $\bar{x} \pm s = 380 \pm 840 \mu\text{g/L}$ ). However, cobwebs alone did not increase TP values, thus it was classified as a weak contaminant. Also,

uniform contamination sources such as dustfall and pollen are classified as none contaminants because the distributions of them are spatially uniform.

Based on the above data classification, about 32 percent of measured data were classified as strongly contaminated data and were eliminated from the data sets. Figure 4 shows the change of frequency curves before and after removing contaminated data by this step. Although most of the high TP concentration data were removed using this method, there were still

unreasonably high TP values remaining. These high TP values are possibly the result of missing contamination flags or invisible contamination sources such as body fluids from insects and animals.

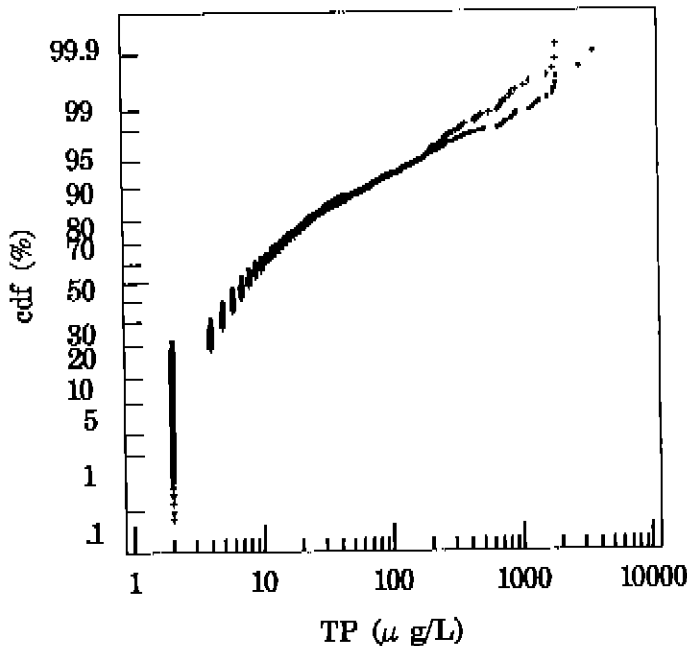


Figure 4. Cumulative Distribution Functions (cdf) of the Pooled Wet TP Concentration Data from 15 Monitoring Sites for Before (•) and After (o) Removing the Contaminated Data by the Step I.

## STEP II: STATISTICAL OUTLIER DETECTION

To identify the outliers in a wet TP data set from a given site, the TP data from nearby sites were taken as the explanatory (independent) variable in Equation (1). Before applying the HS method, the dependent variable,  $Y = \{y_i, i = 1, \dots, N\}$ , at a given site was divided into two subsets: a complete data set  $Y_C = \{y_i, i = 1, \dots, N_C\}$  for which the concurrent measurement at explanatory site is available, or an incomplete data set  $Y_I = \{y_i, i = 1, \dots, N_I\}$  for which  $x_i$  is missing, where  $N (= N_C + N_I)$  is the sample size of  $Y$ . Then, a simple linear regression model of

$$y_i = a + bx_i \quad \text{for } i = 1, \dots, N_C \quad (3)$$

was applied to the complete data set where  $a$  and  $b$  are the regression parameters.

The number of complete data,  $N_C$ , was commonly less than  $N$  due to the large number of missing data that occur randomly. To increase the size of  $N_C$ , one to four explanatory sites were chosen based on the

distance between sites, periods of record, and the number of missing data (Table 2). An average of the TP values measured concurrently from the selected explanatory sites was then taken as an independent variable. For outlier detection modeling, the sites having few high TP values were initially selected as explanatory sites in order to avoid possible outliers in dependent data set. The site ID numbers given in Table 2 are the actual sequence of modeling. In particular, BG2 site was chosen as the explanatory site of the first model because this site has only one obvious outlying value that exceeds about 6000  $\mu\text{g/L}$ .

For each site, the IIS method was first applied to the complete data set  $Y_C$ , from which outliers and a site-specific cutoff value were determined. The site specific cutoff value was then used to identify the outliers in an incomplete data set  $Y_I$ . Table 2 presents the summary statistics of culled data, where the outlier bound given in the fourth column is the largest inlier in  $Y_C$  after removing outliers.

Both  $x_i$  and  $y_i$  were log-transformed (natural logarithm) to help meet a normality condition because the data were positively skewed in most cases (be discussed later). In particular, the chi-squared statistics ( $\chi^2$ ) of both with and without log-transformed clean data were computed for normality test. The result of  $\chi^2$  test in Table 2 reveals that 12 sites accept the hypothesis that the log-transformed data are normally distributed at 95 percent probability level, and that the log-transformed data of 14 sites (except BG1) pass the normality test at 99.5 percent probability level. On the other hand, for the non-transformed data, only three sites satisfy a normality assumption of the data, justifying the log-transformation of the data in Step II.

To assess the performance of the HS (1993) method, Figure 5 plots the estimated  $d_i$  versus TP value for four arbitrarily selected sites (although all sites showed nearly identical patterns). Since the marks in each plot are for the complete data having a size of  $N_C$ , the number of outliers displayed in each plot is slightly less than the corresponding number listed in the third column in Table 2. For the BG1 site, there was no outlier by the HS method, while one observation exceeded the limiting value by the conventional method. For the L-6 and remaining sites, the results by both methods are the same. In particular, the plot for the BG2 site shows three different probability levels (10 percent, 1 percent, 0.1 percent) for investigating the sensitivity of the results of outlier detection to the t-statistics.

As shown in this plot, the result of outlier detection is not sensitive to the significance level,  $\alpha$ , as was also true at the other sites. This insensitivity is due to the fact that most outliers are substantially larger than inliers. Also, due to such distinct differences, the



TABLE 2. Summary of Outlier Detection for Wet TP Concentration Data, Where  $n_0$  is the Numbers of Measured Samples,  $n_{II}$  is the Number of Outliers Identified by the Step II, and  $n$  is the Number of Clean Samples.  
 +: The tabulated  $\chi^2$  value with 2 degrees of freedom.

Site ID and Name for Y	Site IDs for X	$n_0$	$n_{II}$	$n$	Outlier Bound	Mean ( $\mu\text{g/L}$ )	S.D. ( $\mu\text{g/L}$ )	$\chi^2$ Value	
								TP	ln(TP)
1. BG1	2	111	1	86	> 74	8.4	11.1	114.2	12.0
2. BG2	1	114	4	85	> 117	11.0	15.2	89.1	4.8
3. ENR	4,5,6	170	4	69	> 141	12.0	16.6	54.3	10.1
4. OKEEFS	3,5,6	110	2	75	> 40	7.5	7.1	66.9	1.5
5. S-140	3,4,6	88	0	62	> 66	8.1	9.9	57.3	2.6
6. S-7	3,5,8	137	1	64	> 51	9.0	9.1	1.3	1.3
7. S-65A	3,4,5,6	143	9	84	> 164	16.8	22.9	5.1	3.3
8. ENPRC	1,2,3,5	92	0	57	> 35	8.2	10.9	53.3	9.5
9. G-96	1,2,3	68	3	51	> 88	16.6	21.2	72.3	3.2
10. S-127	4,9	98	5	66	> 124	23.9	33.8	81.3	3.7
11. S-131	1,2,9,10	80	3	64	> 175	13.0	23.6	54.9	5.4
12. S-310	1,2,10,11	86	1	74	> 94	11.6	17.5	6.0	2.6
13. L-67A	3,5,6,8	23	0	15	> 14	5.6	3.9	88.8	2.6
14. L-6	5,6,13	33	1	17	> 28	7.9	6.9	104.9	9.9
15. S-308	1,3,12	66	7	51	> 165	18.0	17.7	47.1	2.4
Sum/ Average		1419	41	920	> 92	11.8	15.2	$\chi^2_{.95} = 6.0$ or $\chi^2_{.995} = 10.1$	

results of outlier detection were consistent regardless of the sample size and the number of detected outliers (although the  $R^2$  in regression was low in some cases;  $R^2 = 0.35\sim 0.78$ ). In other words, the goodness-of-the-fit of a regression is not critical for outlier detection.

After removing the identified outliers by this step (about 3 percent of observed data), the sample mean and standard deviation of the wet TP data from 15 sites were estimated to be 11.8  $\mu\text{g/L}$  and 15.2  $\mu\text{g/L}$ , respectively, while those for the data before the removal (and after Step I analysis) were 30.4  $\mu\text{g/L}$  and 92.3  $\mu\text{g/L}$ , respectively.

#### A LUMPED CUTOFF VALUE

As noted in the previous section, the cutoff values for determining outliers varied from site to site depending on the occurrence and magnitude of data contaminations which were, by and large, random in space and time. However, it may be useful to know a lumped cutoff value that can be applicable to all sites. The lumped cutoff value, for instance, can be used to detect outliers at the remaining four sites (ENR-101, ENR-203, ENR-301, ENR-401) or as a screening or

cautioning indicator of future wet TP concentration data, avoiding the need for further statistical analyses.

To identify a lumped cutoff value, a trial and error method was used. That is, with several assumed cutoff values ranging from 100  $\mu\text{g/L}$  to 150  $\mu\text{g/L}$ , the data with less than a given cutoff value were taken from all 15 sites, from which the mean and standard deviation were computed respectively by:

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m \bar{x}_i = \frac{1}{m} \sum_{i=1}^m \left[ \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij} \right] \quad (4)$$

$$s = \frac{1}{m} \sum_{i=1}^m \left[ \frac{1}{(n_i - 1)} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 \right]^{1/2} \quad (5)$$

where  $m (=15)$  and  $n_i$  are the number of sites and the number of samples at site  $i$ , respectively, and  $x_{ij}$  denotes the  $j$ -th wet TP concentration value at site  $i$ . An optimal cutoff value of 130  $\mu\text{g/L}$  was then determined (Figure 6) by matching the computed statistics (by Equations 4 and 5) with the corresponding average statistics of the screened data obtained from the

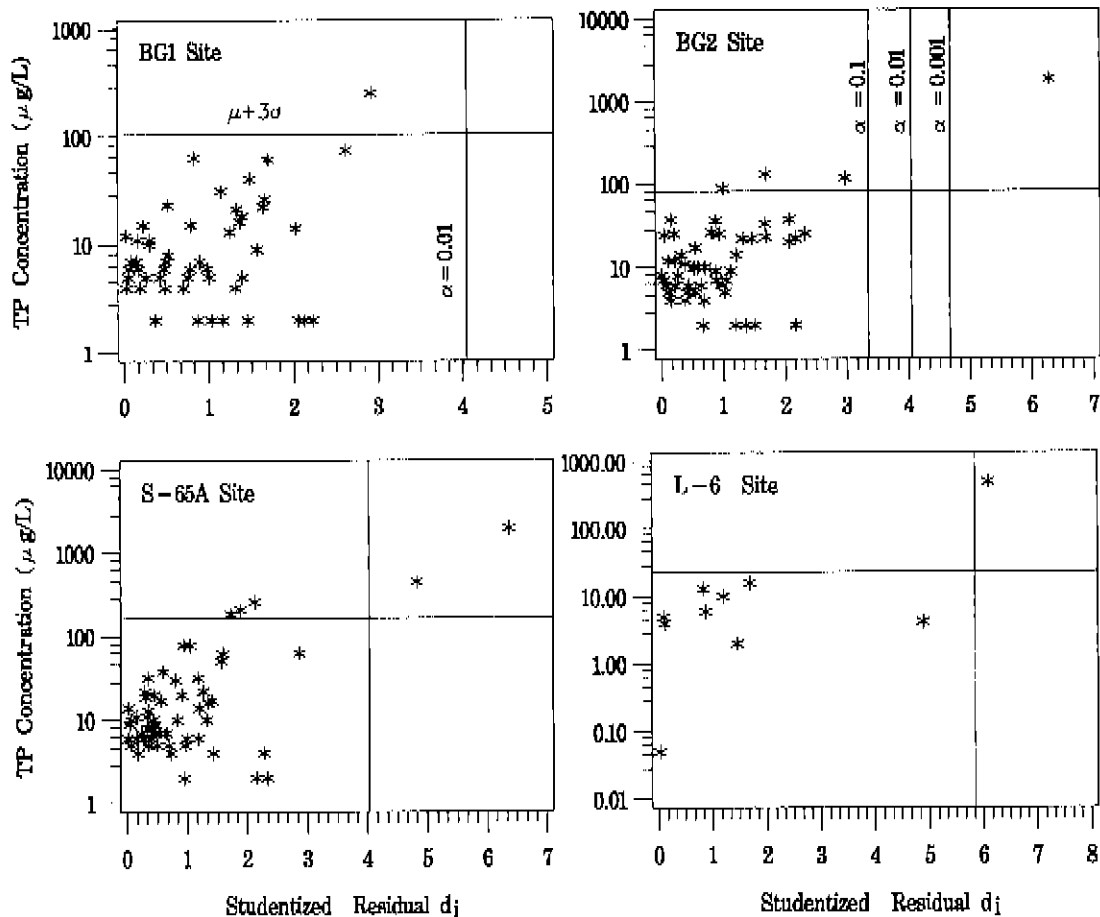


Figure 5. Examples of Outlier Detection for the Selected Sites, Where the Vertical Reference Line is the t-Value for  $\alpha = 0.01$  Unless Noted, While the Horizontal Reference Line is a Simple Conventional Outlier Detection Method Based on the  $\mu + 3\sigma$  criterion.

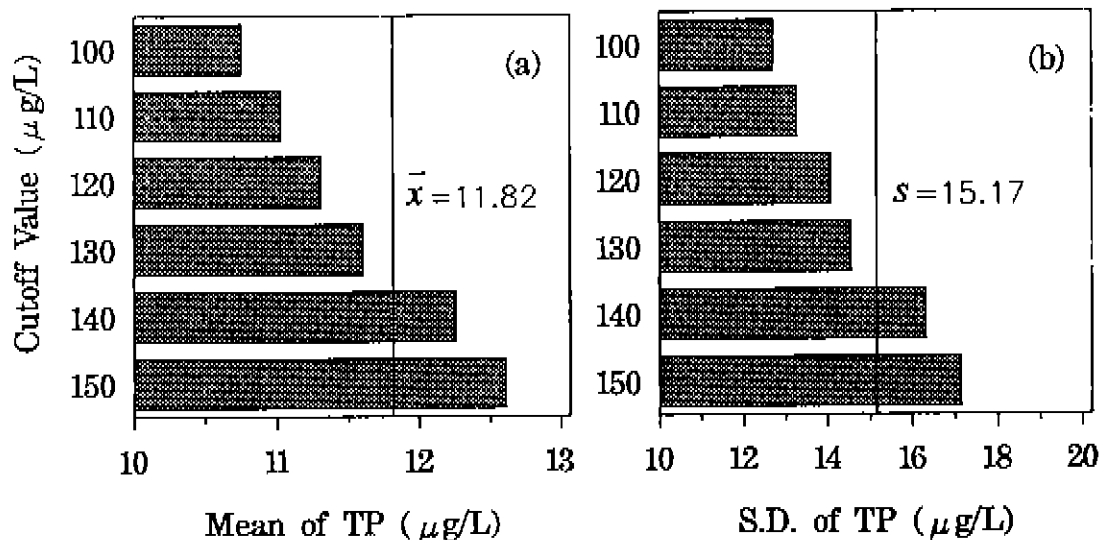


Figure 6. Sample Statistics of Wet TP Data from 16 Atmospheric Deposition Sites Based on Several Alternative Cutoff Values to Determine a Lumped Cutoff Value.

distributed cutoff values given in Table 2. In contrast to the lumped cutoff value, the average cutoff value from 15 sites (the fourth column in Table 2) is 92  $\mu\text{g/L}$ . However, this value is less reliable because it is based on the largest inlier.

## DISCUSSION

An accurate account of phosphorus loads is needed to understand its impact on the Everglades in South Florida (Davis, 1994; Redfield, 1998). This account must include surface loads and atmospheric deposition. The atmospheric deposition component is significant but contamination from bird droppings, body parts of insects, and miscellaneous debris, as documented here, is problematic in obtaining accurate background rates of phosphorus deposition. The method employed here is useful and defensible in removing the bias of contamination for multi-site data.

The estimate of wet TP concentration,  $11.8 \pm 15.2$   $\mu\text{g/L}$ , is consistent with estimates from the Loxahatchee National Wildlife Refuge in Florida of 14  $\mu\text{g/L}$  (Walker and Jewell, 1997), and from the Florida Atmospheric Mercury Study project of 3–7  $\mu\text{g/L}$  (Landing, 1997). But it is less than the  $52 \pm 89$   $\mu\text{g/L}$  determined in south of Lake Okeechobee in May to June of 1992 (Peters and Reese, 1995). Wet TP deposition load in south Florida, with a mean rainfall of 1.35 m/year (Sculley, 1986) and a mean wet TP concentration of 11.8  $\mu\text{g/L}$ , is estimated as 15.9  $\text{mg/m}^2/\text{year}$ . This load estimate matches the estimate from wet/dry collectors throughout Florida of 11 (6–16)  $\text{mg/m}^2/\text{year}$  (Hendry *et al.*, 1981). These comparisons provide a certain level of confidence regarding the District's sampling network, procedures, and the statistical approach that we have taken. However there is still a large amount of variability within our own data.

Wet deposition is quite variable both in space (Hicks *et al.*, 1993; van Ek and Draaijers, 1994; Dixon *et al.*, 1996; Hendry *et al.*, 1981), and time (Hicks *et al.*, 1993). The latter is primarily a result of episodic events. The spatial and temporal variabilities are also presented in the data from the District's monitoring stations. The standard deviation of the samples is equivalent to the mean (after data screening). Also the means ranged from an average of 5.6  $\mu\text{g/L}$  at south of Water Conservation Areas (L67A) to 23.6  $\mu\text{g/L}$  at S127 which is a site north of Lake Okeechobee.

A major question regarding wet deposition is how much is from external or background atmospheric deposition, and how much is from internal or local sources. This question not only deals with the problem of contamination but also methods of collecting wet deposition. Because of the relative simplicity and robustness of this outlier identification technique, it should be useful for any wet deposition collection method.

## SUMMARY

This study attempted to detect outliers in the measured rainfall-borne phosphorus concentration data in which outliers are very common due to contamination from bird droppings, insects and animals, and miscellaneous debris. The approach used both field notes describing the visual inspection of the samples and outlier detection statistics based on a linear regression model. In particular, the study demonstrated how a two-step outlier detection approach can be applied for multi-site environmental data. The forward search algorithm proposed by Hadi and Simonoff (1993) for finding a clean subset was fast and robust as was reported in the previous studies. It was also found through this study that this method was not sensitive to the significance level in the outlier detection method. Although this approach cannot remove all uncertainty from these data, it can be a tool for detecting outliers in the wet TP data observed in South Florida.

As a result of data screening, about 35 percent of the observed data were identified as contaminated and were removed for further data analyses. The averaged mean and standard deviation of the wet TP data collected from 15 sites (after removing the outliers) was 11.8  $\mu\text{g/L}$  and 15.2  $\mu\text{g/L}$ , respectively. Also identified in this study was a lumped cutoff value of 130  $\mu\text{g/L}$ . This lumped cutoff value may be useful for detecting outliers at other sites and for the quality control of future atmospheric deposition sampling.

## ACKNOWLEDGMENTS

The author is grateful to Cheol Mo and Maria Loucraft-Manzano for discussions early in the work; Garth Redfield, Thomas James, Susan Gray, and Linda Lindstrom from the District for their valuable comments through internal peer reviews on the draft manuscript; and three anonymous referees for their constructive comments and corrections.

## LITERATURE CITED

- Asman, W. A. H., T. B. Ridder, H. R. Reijnders, and J. Slanina, 1982. Influence and Prevention of Bird-Droppings in Precipitation Chemistry Experiments. *Water, Air and Soil Pollution* 17:415-420.
- Atkinson, A. C., 1994. Fast Very Robust Methods for the Detection of Multiple Outliers. *Journal of the American Statistical Association* 89:1329-1339.
- Atkinson, A. C. and H. M., Mulira, 1993. The Statistical Plot for the Detection of Multivariate Outliers. *Statistics and Computing* 3:27-35.
- Barnett, V. and T. Lewis, 1984. *Outliers in Statistical Data* (Second Edition). John Wiley, New York, New York.
- Beckman, R. J. and R. D. Cook, 1983. Outliers .....a. *Technometrics* 25(2):119-149.
- Bigelow, D. S., 1984. *Instruction Manual: NADP/NTN Site Selection and Installation*. Natural Resource Ecology Laboratory, Colorado State University, Fort Collins, Colorado, pg. 23.
- Bigelow, D. S. and S. R. Dossett, 1988. *NADP/NTN Instruction Manual: Site Operation*. Natural Resource Ecology Laboratory, Colorado State University, Fort Collins, Colorado, pg. 114.
- Chib, S. and R. C. Tiwari, 1994. Outlier Detection in the State Space Model. *Statistics and Probability Letters* 20:143-148.
- Cook, R. D., 1977. Detection of Influential Observation in Linear Regression. *Technometrics* 19(1):15-18.
- Davis, S. M., 1994. Phosphorus Inputs and Vegetation Sensitivity in the Everglades *In: Everglades - the Ecosystem and Its Restoration*, S. M. Davis and J. C. Ogden (Editors). St. Lucie Press, St. Lucie, Florida, pp. 357-378.
- Dixon, L. K., S. Murray, J. S. Perry, P. J. Minotti, M. S. Henry, and R. H. Pierce, 1996. *Assessment of Bulk Atmospheric Deposition to the Tampa Bay Watershed*. Final Report submitted to the Tampa Bay National Estuary Program, St. Petersburg, Florida.
- Grubbs, F. E., 1969. Procedures for Detecting Outlying Observations in Samples. *Technometrics* 11:1-21.
- Hadi, A. S., 1992. Identifying Multiple Outliers in Multivariate Data. *J. R. Statist. Soc. B*, 54:761-771.
- Hadi, A. S., 1994. A Modification of a Method for the Detection of Outliers in Multivariate Samples. *J. R. Statist. Soc. B*, 56(2): 393-396.
- Hadi, A. S. and J. S. Simonoff, 1993. Procedures for the Identification of Multiple Outliers in Linear Models. *Journal of the American Statistical Association* 88:1264-1272.
- Hadi, A. S. and M. S. Son, 1990. Some Properties of, and Relationships Among, Several Uncorrelated and Homoscedastic Residual Vectors. *Communications in Statistics, Part A - Theory and Methods*, Vol. 19, pp. 2625-2642.
- Hawkins, D. M. and J. S. Simonoff, 1993. High Breakdown Regression and Multivariate Estimation. *Applied Statistics* 42:423-432.
- Hendry, C.D., P. L. Brezonik, and E. S. Edgerton, 1981. Atmospheric Deposition of Nitrogen and Phosphorus in Florida. *In: Atmospheric Pollutants in Natural Waters*, S. J. Eisenreich (Editor). Ann Arbor Sci., Publ., Ann Arbor, Michigan, pp. 199-215.
- Hicks, B., R. McMillan, R. S. Turner, G. R. Holdren Jr., and T. C. Strickland, 1999. A National Critical Loads Framework for Atmospheric Deposition Effects Assessment: III. Deposition Characterization. *Environmental Management* 17:343-353.
- Jain, R. B., 1981. Percentage Points of Many - Outliers Detection Procedures. *Technometrics* 23(1):71-75.
- Landing, W. M., 1997. Measurement of Aerosol Phosphorus in South Florida. *In: Proceedings of a Conference on Atmospheric Deposition into South Florida*, October 1997. South Florida Water Management District (SFWMD), West Palm Beach, Florida.
- Peters, N. E. and R. S. Reese, 1995. Variations of Weekly Atmospheric Deposition for Multiple Collectors at a Site on the Shore of Lake Okeechobee, Florida. *Atmospheric Environment* 29:179-187.
- Redfield, G. W., 1998. *Quantifying Atmospheric Deposition of Phosphorus: A Conceptual Model and Literature Review for Environmental Management*. Technical Publication WRE No. 360, SFWMD, West Palm Beach, Florida, 35 pp.
- Rousseeuw, P. J. and B. C. Van Zomeren, 1990. Unmasking Multivariate Outliers and Leverage Points (with discussion). *Journal of the American Statistical Association* 85:633-651.
- Sculley, S., 1986. *Frequency Analysis of SFWMD Rainfall*. Technical Publication 86-6, SFWMD, West Palm Beach, Florida.
- SFWMD, 1996. *SFWMD Comprehensive Quality Assurance Plan*, WRE-346, SFWMD, West Palm Beach, Florida.
- Tiwari, R. C. and T. P. Dienes, 1994. The Kalman Filter Model and Bayesian Outlier Detection for Time Series Analysis of BOD Data. *Ecological Modelling* 73:159-165.
- van Ek, R. and G. P. J. Draaijers, 1994. Estimates of Atmospheric Deposition and Canopy Exchange from Three Common Tree Species in the Netherlands. *Water, Air, and Soil Pollution* 73:61-82.
- van Wyk, D. B. and W. D. Stock, 1991. Design of a Sequential Atmospheric Deposition Sampler for Use in Remote Catchments. *Water SA* 17(3):183-188.
- Walker, W. W. and S. D. Jewell, 1997. Atmospheric Deposition of Phosphorus in Loxahatchee National Wildlife Refuge. *In: Proceedings of a Conference on Atmospheric Deposition into South Florida*. SFWMD, West Palm Beach, Florida.
- Woodruff, D. and D. M., Rocke, 1993. Heuristic Search Algorithms for Minimum Volume Ellipsoids. *J. of Computational and Graphical Statistics* 2:69-95.